

Lecture 1

Name: Tyler LaBonte

Professor: Santosh S. Vempala

1 Statistical Inference

Suppose $x_1, x_2, \dots, x_m \in \mathbb{R}$. Then we can use the following statistics:

- Mean: $\mu = \frac{1}{m} \sum x_i$,
- Variance: $\sigma^2 = \frac{1}{m} \sum (x_i - \mu)^2$.

In high dimension, say \mathbb{R}^d , we must be precise about variance. We will project each point onto a unit vector $\mathbf{v} \in \mathbb{R}^d$:

$$\sigma_{\mathbf{v}}^2 = \frac{1}{m} \sum_i (\mathbf{x}_i \cdot \mathbf{v} - \boldsymbol{\mu} \cdot \mathbf{v})^2 \quad (1)$$

$$= \frac{1}{m} \sum_i ((\mathbf{x}_i - \boldsymbol{\mu}) \cdot \mathbf{v})^2 \quad (2)$$

$$= \mathbf{v}^\top \left(\frac{1}{m} \sum_i (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top \right) \mathbf{v}. \quad (3)$$

The middle term is $\boldsymbol{\Sigma}$, the covariance matrix. We can write

$$\sigma_{\mathbf{v}}^2 = \mathbf{v}^\top \boldsymbol{\Sigma} \mathbf{v}. \quad (4)$$

Of course, the mean and variance does not characterize the data distribution. One way to properly characterize the distribution is by estimating the probability density function (PDF). For example, the uniform distribution in $[a, b]$ has density

$$p(x) = \begin{cases} \frac{1}{|b-a|} & x \in [a, b], \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$

By definition, a PDF is a nonnegative integrable function which integrates to 1 over the domain.

2 The Gaussian Distribution

2.1 The Standard Gaussian

The one-dimensional standard Gaussian $\mathcal{N}(0, 1)$ has PDF

$$p(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}. \quad (6)$$

In \mathbb{R}^d , the distribution is called the multivariate standard Gaussian $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. The PDF is the product of 1-dimensional Gaussians along each coordinate:

$$p(\mathbf{x}) = \frac{1}{(\sqrt{2\pi})^d} e^{-\|\mathbf{x}\|^2/2} \quad (7)$$

$$= \prod_i \frac{1}{\sqrt{2\pi}} e^{-x_i^2/2}. \quad (8)$$

Fact 2.1: Spherical Symmetry of the Multivariate Standard Gaussian

For $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and any unit vector $\mathbf{v} \in \mathbb{R}^d$, we have $\mathbf{x} \cdot \mathbf{v} \sim \mathcal{N}(0, 1)$. Equivalently,

$$\Pr[\mathbf{x} \cdot \mathbf{v} = t] = \frac{1}{\sqrt{2\pi}} e^{-t^2/2}. \quad (9)$$

Proof: We have

$$\Pr[\mathbf{x} \cdot \mathbf{v} = t] = \int_{\mathbf{x}: \mathbf{x} \cdot \mathbf{v} = t} \frac{1}{(\sqrt{2\pi})^d} e^{-\|\mathbf{x}\|^2/2} d\mathbf{x}. \quad (10)$$

We can write \mathbf{x} in the basis $(\mathbf{v}, \mathbf{v}_2, \dots, \mathbf{v}_d)$ so that

$$\|\mathbf{x}\|^2 = (\mathbf{x} \cdot \mathbf{v})^2 + \|\mathbf{y}\|^2. \quad (11)$$

Then,

$$\Pr[\mathbf{x} \cdot \mathbf{v} = t] = \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \int_{\mathbf{y} \in \mathbb{R}^{d-1}} \frac{1}{(\sqrt{2\pi})^{d-1}} e^{-\|\mathbf{y}\|^2/2} d\mathbf{y}. \quad (12)$$

Since the integrand is a $(d-1)$ -dimensional Gaussian, its integral over \mathbb{R}^{d-1} is 1 by definition. We have to be careful about the differential here by checking the determinant of its Jacobian, but since $(\mathbf{v}_2, \mathbf{v}_3, \dots, \mathbf{v}_d)$ is orthogonal, we immediately have $d\mathbf{y} = d\mathbf{x}$. ■

Question 2.1

Where does the $\frac{1}{\sqrt{2\pi}}$ come from in the denominator of the PDF?

Fact 2.2: Gaussian PDF Integrates to 1

Let $p(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$. Then, $\int_{\mathbb{R}} p(x)dx = 1$.

Proof: This is a difficult integral *a priori*, so we consider the two-dimensional case:

$$p(x, y) = \frac{1}{2\pi}e^{-(x^2+y^2)/2} \quad (13)$$

$$\int_{\mathbb{R}^2} p(x, y)dxdy = \left(\int_{\mathbb{R}} p(x)dx \right)^2. \quad (14)$$

We take advantage of spherical symmetry by integrating in spherical coordinates (in two dimensions, these are concentric circles).

$$\left(\int_{\mathbb{R}^2} p(x)dx \right)^2 = \frac{1}{2\pi} \int_{r=0}^{\infty} 2\pi r e^{-r^2/2} dr \quad (15)$$

$$= \int_{r=0}^{\infty} r e^{-r^2/2} dr. \quad (16)$$

Let $u = r^2/2$, then

$$\int_{r=0}^{\infty} r e^{-r^2/2} = \int_{u=0}^{\infty} e^{-u} du \quad (17)$$

$$= 1. \quad (18)$$

■

2.2 The General Gaussian

The one-dimensional Gaussian $\mathcal{N}(\mu, \sigma^2)$ has PDF

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right). \quad (19)$$

The multivariate Gaussian $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ has PDF

$$p(\mathbf{x}) = \frac{1}{(\sqrt{2\pi})^d |\det(\boldsymbol{\Sigma})|} \exp\left(-\frac{(\mathbf{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}{2}\right), \quad (20)$$

with covariance matrix $\boldsymbol{\Sigma}$ positive definite (since we cannot have negative variance).

In a standard Gaussian, the densities are symmetric over the sphere. In a general Gaussian, they are symmetric over an ellipsoid. That is, a general Gaussian is a standard Gaussian acted on by a linear transformation described by $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$.

Fact 2.3: Linear Transformations of a Gaussian

Let $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. Then for $\mathbf{A} \in \mathbb{R}^{d \times d}$ and $\mathbf{b} \in \mathbb{R}^d$, $\mathbf{Ax} + \mathbf{b} \sim \mathcal{N}(\mathbf{b}, \mathbf{AA}^\top)$. In the other direction, let $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Then $\boldsymbol{\Sigma}^{\frac{1}{2}}(\mathbf{x} - \boldsymbol{\mu}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$.

3 Maximum Likelihood Estimation for Gaussians

Suppose $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m \in \mathbb{R}^d$ and we would like to find the best-fit Gaussian. In the case when $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}$, we would like to find $\mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$ which maximizes the likelihood of generating the observed data. The log-likelihood function we will maximize is

$$L(\boldsymbol{\mu}, \sigma^2) = \log \prod_{j=1}^m \frac{1}{(\sqrt{2\pi}\sigma)^d} \exp\left(\frac{-\|\mathbf{x}_j - \boldsymbol{\mu}\|^2}{2\sigma^2}\right) \quad (21)$$

$$= -\left(md \log(\sqrt{2\pi}) + md \log(\sigma) + \frac{1}{2\sigma^2} \sum_j \|\mathbf{x}_j - \boldsymbol{\mu}\|^2\right). \quad (22)$$

Lemma 3.1: Maximum Likelihood Estimation for Gaussians

The sample mean $\hat{\boldsymbol{\mu}}$ and sample variance $\hat{\sigma}^2$ are maximum likelihood estimators for the population mean $\boldsymbol{\mu}$ and population variance σ^2 .

Proof: For the variance,

$$\frac{\partial L}{\partial \sigma} = -\frac{md}{\sigma} + \frac{1}{\sigma^3} \sum_j \|\mathbf{x}_j - \boldsymbol{\mu}\|^2. \quad (23)$$

Setting equal to zero implies

$$\sigma^2 = \frac{1}{md} \sum_j \|\mathbf{x}_j - \boldsymbol{\mu}\|^2 = \hat{\sigma}^2. \quad (24)$$

For the mean,

$$\frac{\partial L}{\partial \mu_i} = -\frac{1}{2\sigma^2} \left(-2 \sum_j (x_{ji} - \mu_i) \right). \quad (25)$$

Setting equal to zero implies

$$\mu_i = \frac{1}{m} \sum_j x_{ji} = \hat{\mu}_i. \quad (26)$$

■

4 The Central Limit Theorem

The Central Limit Theorem (CLT) is a fundamental result which states that the distribution of sums (or equivalently, means) of independent random variables tends towards a Gaussian. There

are many quantitative formulations of the CLT, but the following provides a useful bound based on the first three moments.

Theorem 4.1: Berry-Esseen CLT

Let X_1, X_2, \dots, X_n be independent random variables. Let $Y_n = \sum_{i=1}^n X_i$ and $Z_n \sim \mathcal{N}(\mathbb{E}[Y_n], \text{Var}(Y_n))$. Then for $t \in \mathbb{R}$ and $c \in [0, 1)$,

$$|\Pr[Y_n \leq t] - \Pr[Z_n \leq t]| \leq \frac{c}{\text{Var}(Y_n)^{3/2}} \sum_i \mathbb{E}[|X_i|^3]. \quad (27)$$

Example. Suppose $X_1, X_2, \dots, X_n \sim \text{Ber}(-1, 1)$. Then $\mathbb{E}[|X_i|^3] = 1$ and $\text{Var}(Y_n) = n$. By the theorem,

$$|\Pr[Y_n \leq t] - \Pr[Z_n \leq t]| \leq \frac{c}{n^{3/2}} n \quad (28)$$

$$= \frac{c}{\sqrt{n}}. \quad (29)$$